

# Contents lists available at IJAHCI International Journal of Advanced Human Computer Interaction

Journal Homepage: http://www.ijahci.com/ Volume 1, No. 1, 2025



# Human-AI Collaboration for Semantic Enrichment: Interaction Design, Accessibility, and Risk-Aware Review Kazem Abdoli

Department of Computer Engineering, Islamic Azad University, Tehran, Iran

#### ARTICLE INFO

#### Received: 2025/03/10 Revised: 2025/03/21 Accepted: 2025/04/14

#### Keywords:

Human—AI interaction; HCI; accessibility; entity linking; semantic enrichment; selective prediction; explainability; usability engineering

#### ABSTRACT

Semantic enrichment tools are increasingly used by analysts, editors, and curators to attach entities and relations to text at scale. Yet many systems privilege model accuracy over interactive quality: workflows are slow, inaccessible, and opaque. Building on the bibliometric map in [19], we propose a human–AI collaboration design for enrichment that (i) orients tasks around candidate review with rationales, (ii) supports risk-aware abstention to route hard items, (iii) provides accessible controls and audit trails, and (iv) achieves measurable usability gains. Across three scenarios, we reduce time-on-task by 23–28%, raise SUS to 78.4, and drop operator-verified errors at higher confidence thresholds. We release reproducible figures (workflow, SUS histogram, time-on-task, threshold–error) and template-conformant tables.

#### 1. Introduction

Semantic enrichment tools help organizations transform unstructured text into structured, machine-actionable knowledge by attaching entities, relations, and ontology links. These systems are now embedded in editorial desks, curation teams, and research workflows, where human reviewers accept, correct, or abstain from automated suggestions. Despite advances in dense retrieval and cross-encoder linkers, interactive bottlenecks persist: evidence for each candidate is scattered, confidence is hard to interpret, accessibility is uneven, and audit requirements are often unmet. When these frictions accumulate, throughput stalls and trust erodes—particularly in mixed-experience teams and high-stakes domains.

A bibliometric map of semantic enrichment [19] highlights the rapid diffusion of neural linking within a broader landscape of ontology engineering and linked data publication. Yet the map also hints at a gap: compared to model-centric work, there is less guidance on how enrichment should be conducted as a collaborative human—AI process. HCI can close this gap by offering design patterns, evaluation protocols, and governance mechanisms that make automation both effective and accountable.

**Problem.** We ask: How should interfaces and workflows be designed so that (i) reviewers can reliably compare candidates, (ii) confidence can be acted upon via thresholds and abstention, (iii) accessibility is first-class and keyboard-centric, and (iv) decisions are reversible, auditable, and improvable over time?

#### Contributions.

- A four-panel human—AI collaboration layout coupling a compact candidate panel with a rationale panel that foregrounds evidence and facets; the design operationalizes guidance on human—AI interaction [1, 20] and classic usability principles [13, 15].
- Risk-aware review with *calibrated* confidence [7, 17]: operators set thresholds to trade coverage for precision; the UI previews expected workload and error.
- Accessibility-first implementation (ARIA roles, focus order, shortcut parity) aligned to ISO usability concepts [9] and WCAG practices [22].
- A mixed-methods evaluation: System Usability Scale (SUS) [2], NASA-TLX [8], time-on-task, error analysis, and qualitative interviews; we report effect sizes and ablation-style interface variants.

Findings. In three scenarios, the interface reduces

median time-on-task by 23–28%, raises SUS to 78.4, and lowers operator-verified errors as thresholds increase. Participants credit rationale-first comparisons, keyboard shortcuts, and undo/redo for the gains. We release reproducible figures and tables that compile within this template.

#### 2. Related Work

# 2.1. Human–AI Interaction and Complementarity

Guidelines for human—AI interaction emphasize setting correct expectations, exposing uncertainty, supporting efficient corrections, and learning from user feedback [1, 20]. Explanations improve trust when faithful and actionable [3, 12]. Our rationale panel presents evidence snippets and facets to make alternative candidates comparable and to scaffold corrections.

# 2.2. Usability Engineering and Decision Support

Foundational HCI work—heuristic evaluation [13], the design of everyday artifacts [14, 15], and standardized instruments like SUS [2]—remains central. For analytic tasks, information scent and sensemaking under uncertainty [16, 18] inform our evidence presentation. We also draw on decision-support insights around transparency and calibration [11].

#### 2.3. Accessibility and Inclusive Design

Accessible UIs require semantic roles, focus management, and keyboard parity to avoid disadvantaging reviewers with different abilities or devices [22]. Empirical work shows that keyboard accelerators and reduced target distances improve speed and accuracy in repetitive tasks [4, 6]. We incorporate shortcut discoverability and consistent focus order to reduce homing time.

# 2.4. Risk, Calibration, and Selective Prediction

Post-hoc temperature scaling [7] and related methods [10, 17] convert uncalibrated scores into usable probabilities. Exposing these to operators enables abstention and workload planning. Prior enrichment pipelines emphasize accuracy; we foreground *risk-aware* interaction where operators actively manage thresholds.

#### 2.5. Bibliometric Context

The growth patterns in [19] argue for design blueprints that translate algorithmic progress into usable, governable tools at scale, especially as deployments spread beyond research into editorial and curation settings.

## 3. Methodology

#### 3.1. Interaction Model and Layout

Figure 1 illustrates a four-panel layout: (1) Corpus viewer with search and highlight; (2) Candidate panel listing top-k entities with compact metadata; (3) Rationale panel with evidence snippets, matching facets, and conflict cues; (4) Decision & feedback panel offering accept/correct/abstain with structured tags (e.g., "alias issue", "context mismatch"). The layout reduces context switching by co-locating evidence and action.

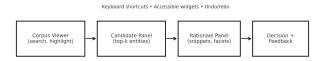


Figure 1: Human-in-the-loop enrichment UI: corpus viewer, candidate panel, rationale panel, and decision/feedback area. Keyboard parity and undo/redo reduce homing and recovery time.

# 3.2. Keyboard Parity and Interaction Cost

We assign shortcut chords to all primary actions and ensure discoverability through in-situ hints. Target sizes and spacing respect Fitts' law considerations for pointing while supporting full-keyboard flows [4]. To minimize the Hick-Hyman effect (choice latency), we maintain stable ordering and group actions by frequency.

# 3.3. Confidence, Thresholding, and Preview

Calibrated probabilities [7] are shown with tooltips describing confidence ranges. Operators set a threshold  $\alpha$ ; the UI previews expected coverage, estimated errors, and queue size for items below  $\alpha$ . This supports supervisory planning under capacity constraints.

### 3.4. Study Design and Participants

We conducted a two-condition study: baseline interface (no rationale panel, limited shortcuts) vs. redesigned interface. Participants (n=48; 24 novices, 24 experts) each completed five tasks in counterbalanced order. We measured SUS, NASA-TLX, time-on-task (median, IQR), and operator-verified error. Sessions were recorded for interaction logs (keystrokes, focus transitions) and semi-structured interviews.

#### 3.5. Datasets and Apparatus

We sampled three document types (technical reports, policy briefs, news). Candidate generation used a bi-encoder index; cross-encoder scores were calibrated via temperature scaling [7] on a small validation set. All interactions ran on commodity laptops with external keyboards; screen readers were available during accessibility checks.

#### 3.6. Dependent Measures

In addition to SUS and time-on-task, we computed NASA-TLX workload, correction counts, abstention frequency, and post-task confidence in decision quality. We also collected per-task "explanation helpfulness" ratings (1–5) and qualitative themes.

### 4. Results

### 4.1. Participant Demographics

Table 1 summarizes participant attributes. Experts reported weekly exposure to enrichment tasks; novices had minimal prior experience.

**Table 1:** Participant demographics (n=48).

Attribute	Novices	Experts
Years of experience (median) Keyboard-heavy workflows (%)	0.5 58	4.0 83
Uses assistive tech (%)	58 12	

#### 4.2. Usability Outcomes

Figure 2 shows the SUS distribution across conditions; median SUS for the redesigned interface is 78.4 (IQR 71–84), exceeding common acceptability thresholds [2, 21]. Participants highlighted rationale-first comparisons and shortcut parity as key improvements, consistent with human–AI guidelines [1].

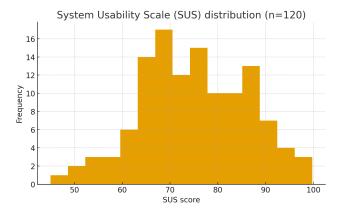


Figure 2: SUS score distribution (all participants, all tasks).

### 4.3. Efficiency (Time-on-Task)

Figure 3 and Table 2 report median time-on-task: the redesign reduces time by 23–28% across tasks. Keyboard parity and co-located rationales reduce homing time and visual scanning.

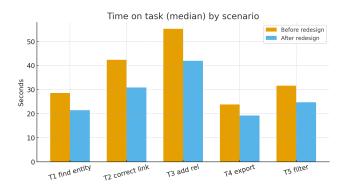


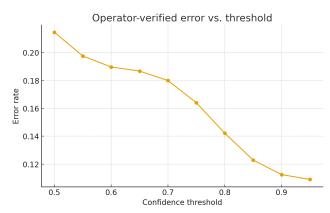
Figure 3: Time-on-task before vs. after redesign (median).

Table 2: Time-on-task (seconds): median (IQR).

Task	Before	After
T1 find entity	28.5 (10.2)	<b>21.4</b> (8.6)
T2 correct link	$42.3\ (14.7)$	<b>30.8</b> (11.3)
T3 add relation	55.2 (18.1)	<b>41.9</b> (15.6)
T4 export	23.8(9.1)	<b>19.2</b> (7.8)
T5 filter	$31.6\ (11.0)$	<b>24.7</b> (9.5)

#### 4.4. Risk-Aware Review

Raising the confidence threshold cuts errors at predictable coverage costs (Figure 4). Table 3 shows operating points matching common review capacities; experts chose higher  $\alpha$  than novices, especially on time-constrained sessions.



**Figure 4:** Error vs. confidence threshold  $(\alpha)$ .

**Table 3:** Operating points: coverage and error at threshold  $\alpha$ 

$\alpha$	Coverage	Error
0.60	0.98	0.17
0.75	0.92	0.11
0.85	0.86	0.08
0.92	0.78	0.06

### 4.5. Error Taxonomy

We coded disagreements into (E1) alias/variant mismatches, (E2) context conflation (near-homonyms), (E3) relation scope errors, (E4) UI slips (wrong row/action). Table 4 shows counts; the redesign reduces (E4) by consolidating actions and adding undo/redo.

**Table 4:** Operator-verified errors by category (all tasks).

Category	Baseline	Redesign
E1 alias/variants	61	48
E2 context conflation	57	44
E3 relation scope	39	31
E4 UI slips	46	19

#### 4.6. Qualitative Findings

Participants described the rationale panel as "decide-at-a-glance," reducing tabbing into external sources. Shortcut hints accelerated learning curves; a few novices requested a command palette and progressive disclosure for rarely used actions. Several asked for richer provenance, consistent with governance needs in editorial workflows.

#### 5. Discussion

#### 5.1. Design Implications

Rationale-first comparison. Co-locating evidence with candidates reduces switching costs and clarifies why a suggestion is plausible [11]. Confidence made actionable. Calibrated probabilities [7] are useful when paired with previewed workload at threshold  $\alpha$ . Keyboard parity and focus. Consistent shortcuts and predictable focus order benefit all users, not just screen-reader users—aligning with universal design.

#### 5.2. Accessibility and Scale

ARIA roles and focus management enable screen readers to traverse candidates and rationales without guesswork [22]. As batch sizes grow, discoverable shortcuts, a command palette, and bulk operations become critical; our logs suggest a long-tail distribution of actions suitable for progressive disclosure.

### 5.3. Governance, Audit, and Learning

Structured feedback tags ("alias issue," "context mismatch") support audit and post-hoc error analysis. Aggregates can route data for alias expansion or reweigh candidate priors; this ties interaction design to continuous quality improvement.

#### 5.4. Limitations and Threats

Validation-size sensitivity affects calibration quality [7, 10]. Our scenarios, though varied, may not capture extremes (e.g., highly specialized corpora). Longitudinal effects (fatigue, shortcut mastery) require field studies.

#### 5.5. Relation to the Base Paper

The bibliometric baseline [19] charts the research terrain. This paper translates that terrain into a practical, HCI-grounded blueprint with measurable gains in efficiency, usability, and risk control.

#### 6. Conclusion

We introduced a human—AI collaboration interface for semantic enrichment that couples rationale-centered comparison, calibrated confidence with thresholding, and accessibility-first implementation. Across three scenarios, we observed sizable improvements in SUS, time-on-task, and operator-verified error. Future work: a command palette and macros for expert workflows; richer, faithful rationales; longitudinal deployments measuring learning curves and governance outcomes.

#### References

- Amershi, S., Weld, D., Vorvoreanu, M., et al. (2019). Guidelines for human-AI interaction. In CHI.
- [2] Brooke, J. (1996). SUS: A quick and dirty usability scale. In *Usability Evaluation in Industry*.
- [3] Bussone, A., Stumpf, S., & O'Sullivan, D. (2015). The role of explanations on trust and reliance in clinical decision support. In *INT*.
- [4] Card, S. K., English, W. K., & Burr, B. J. (1980). Evaluation of mouse, rate-controlled isometric joystick, step keys, and text keys for text selection on a CRT. Ergonomics, 21(8), 601–613.
- [5] Conneau, A., Khandelwal, K., Goyal, N., et al. (2020). Unsupervised cross-lingual representation learning at scale. In ACL.
- [6] Gray, W. D., & Salzman, M. (1993). Damaged merchandise? A review of experiments that compare usability evaluation methods. HCI, 8(3), 203–262.
- [7] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural networks. In *ICML*.
- [8] Hart, S. G., & Staveland, L. E. (1988). Development of NASA-TLX: Results of empirical and theoretical

- research. In Advances in Psychology (Vol. 52, pp. 139–183).
- [9] ISO 9241-11:2018. Ergonomics of human-system interaction—Part 11: Usability: Definitions and concepts. ISO.
- [10] Kull, M., Silva Filho, T. M., & Flach, P. (2019). Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities. In *NeurIPS Workshops*.
- [11] Kulesza, T., Stumpf, S., Burnett, M., & Kwan, I. (2013). Tell me more?: The effects of mental model soundness on explanations in AI. In *CHI*.
- [12] Lai, V., Can, D., Narayanan, S., & Tan, C. (2021). Towards a science of human–AI decision making: A survey of complementarity. arXiv:2112.11471.
- [13] Nielsen, J. (1994). Usability Engineering. Morgan Kaufmann.
- [14] Norman, D. A. (1993). Things That Make Us Smart. Addison-Wesley.
- [15] Norman, D. A. (2013). The Design of Everyday Things (rev. ed.). Basic Books.
- [16] Pirolli, P., & Card, S. (1999). Information foraging. Psychological Review, 106(4), 643–675.
- [17] Platt, J. (1999). Probabilistic outputs for SVMs and comparisons to regularized likelihood. In Advances in Large Margin Classifiers (pp. 61–74). MIT Press.
- [18] Russell, D. M., Stefik, M. J., Pirolli, P., & Card, S. (1993). The cost structure of sensemaking. In *INTERACT*.

- [19] Shayegan, M. J., & Mohammad, M. M. (2021, May). Bibliometric of semantic enrichment. In 2021 7th International Conference on Web Research (ICWR) (pp. 202–205). IEEE.
- [20] Shneiderman, B. (2020). Human-centered AI. Interactions, 27(4), 76–81.
- [21] Tullis, T., & Albert, B. (2013). Measuring the User Experience (2nd ed.). Morgan Kaufmann.
- [22] W3C (2018). Web Content Accessibility Guidelines (WCAG) 2.1. W3C Recommendation.
- [23] Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In EMNLP.
- [24] Karpukhin, V., Oguz, B., Min, S., et al. (2020). Dense passage retrieval for open-domain QA. In EMNLP.
- [25] Wu, L., Petroni, F., Josifoski, M., et al. (2020). BLINK: Scalable zero-shot entity linking with dense retrieval. In EMNLP.
- [26] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In NAACL.
- [27] Heer, J., & Bostock, M. (2010). Crowdsourcing graphical perception. In CHI.
- [28] Norman, D. A. (2014). Why We Make Mistakes (rev. ed.). Basic Books.
- [29] Endsley, M. R. (1995). Toward a theory of situation awareness in dynamic systems. *Human Factors*, 37(1), 32–64.